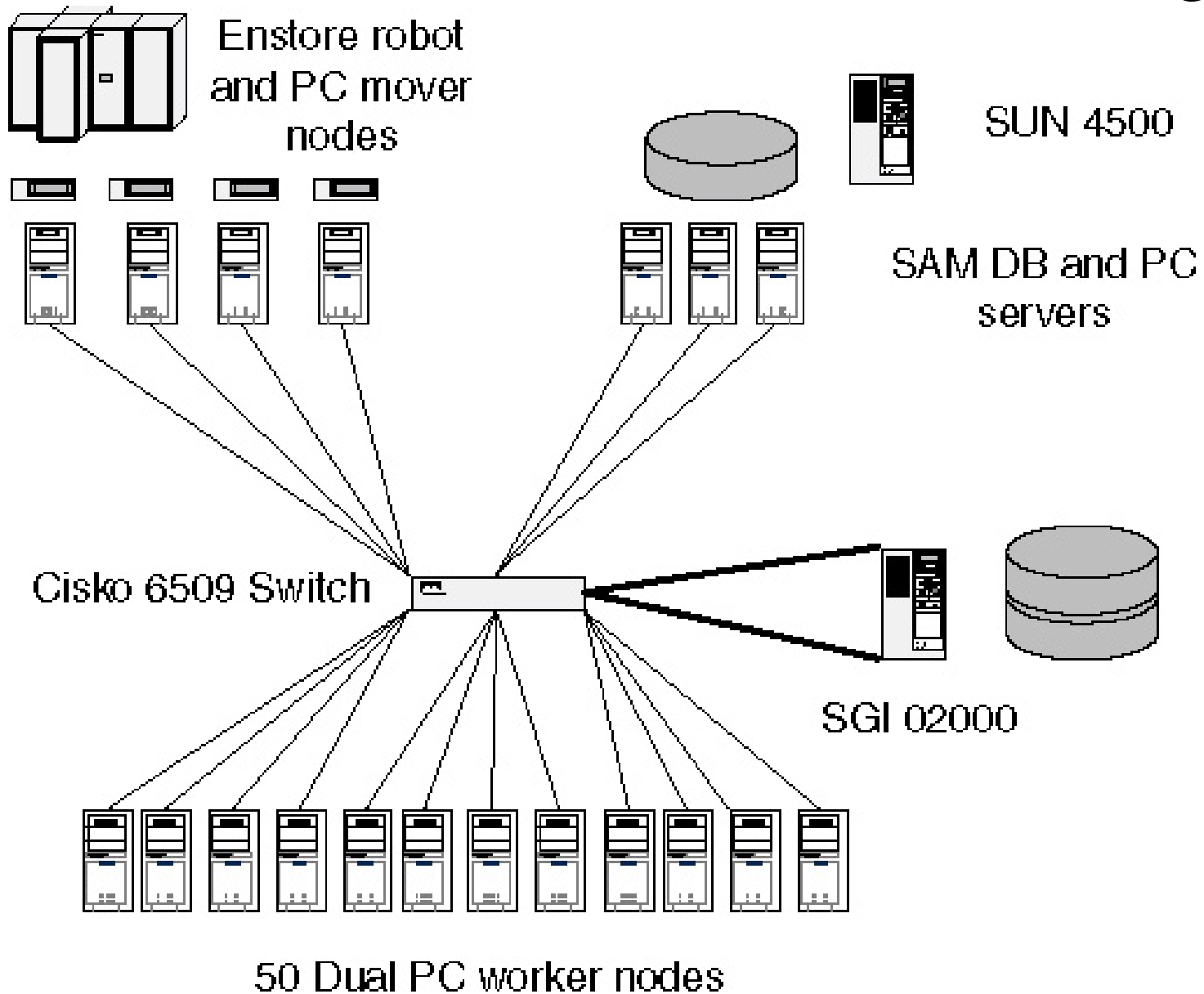# DO Run II Farms

H. Schellman, M. Albert, J. Bakken, L. deBarbaro, M. Breitung,
M. Diesburg, S. Epstein, D. Fagan, J. Fromm, L. Giacchetti, D. Holmgren,
T. Jones, T. Levshina, L. Lueking,  I. Mandrichenko,  S. Mayola,
A. Moibenko,  R. Pasetes, D. Petravick, M.Schweitzer, K. Shepelak,
I. Terekhov, J. Trumbo, S. Veseli,  M. Vranicar,
R. Wellner, S. White, V. White

# DO Farm needs

- ## 250K  event size
- ## 50Hz trigger rate
  - peak rate of  12.5 MB/sec
  - DC is less but reprocessing will bring back up

- ## Reconstruction 5- 10 seconds/event
         ## on 500 MHz PIII
  - need 250-500 CPU's to handle peak rate
  - DC is 40% of peak
  - time constant for  1 GB file is 5- 10 hours.

DO Farms

Enstore robot and PC mover nodes

SUN 4500

SAM DB and PC servers

Cisko 6509 Switch

SGI 02000

50 Dual PC worker nodes

# I/O machine



- **Purpose**
  - split/merge of farm output
  - Serve home areas
  - Batch system control
  - File delivery master
- **D0bbin**
  - 4 CPU SGI 02000
  - 2 GB ethernet cards
  - 4 72 GB disk partitions (2 way stripe)
  - peak I/O rates of 40-60 MB/sec

# Worker Nodes

- Dual Pentium III 500MHz
- 256MB/CPU
- 2 data disks (18 GB) + 6GB system
- Fast ethernet
- CD/floppy for system configuration

# Design Principles

- Use existing facilities
  - SAM/Enstore for data access and file tracking
  - Farm batch system (FBS) for most job control

- Keep D0 farm control scripts to a minimum
  - Batch system assigns machines
  - Data access system decides which file you get

- If worker process or machine dies, lose minimal number of files and don't affect other processes

- No heroic recovery measures, track and resubmit those files

# Worker Configuration

- Workers act as generic FNAL farm machines
  - Only customization is pnfs for file delivery and home area mount
  - D0 environment downloads at job start
  - data access through SAM/encp/rcp, database server

- Batch system assigns workers to job, not D0FARM control process.
- D0FARM control never knows which workers are assigned to a job and does not need to.

# Data Access is SAM/enstore

- Integrated data handling system
- File and process data base
- Data base server
- File servers
- Enstore File delivery systems
- Pnfs file system

Farm Perspective

Can tell it you want a set of files

Can ask for the 'next' file

Can flag file as processed or error

Can get detailed accounting on what happened

**Major Efforts:**

**2-3 talks at this conference**

**SAM # 241**

**Enstore Talk#176**

# Farm Batch System
# Typical Farm Job

SECTION START

    EXEC=startjob  *parameters*

    QUEUE=D0bbin

SECTION WORKER

    EXEC=runjob  *parameters*

    NWORKERS=20

    QUEUE=D0worker

SECTION END

    EXEC=stopjob *parameters*

    QUEUE=D0bbin

    DEPEND WORKER(done)

- Queue tells the system what kind of machine to run on and how many.
- EXEC gives the script name and parameters
- DEPEND allows cleanup section to run when all worker sections are done.
- FBS assigns temporary disk on workers
- On end yanks disk and kills all processes.

# Start Section

- Set up products and output directories on d0bbin
- Tell SAM which files you will want
- Go into wait state until get end signal

# Worker Section

- Download D0 environment
- Start SAM stager
- Ask for next file
- Process file
- Store output file on output buffer
- Inform SAM of success
- Ask for next file
- On error or end of list, terminate.

# End Section

- Create  job summary
- Send message to Start process telling it to shut down the SAM connection for input
- (Optional) Start file merge/store of output files.

## CURRENT JOBS

Status   Monitor   Statistics   Log   Kill   Preference   Print   Save   Quit

| JOB_ID | SECTION_NAME | USER | QUEUE | STATUS | START | FINISHED |
|--------|--------------|------|-------|--------|-------|----------|
| 9150 | *job* | | | | | |
| 9150 | START_SAM | d0farm | io_d0sgi | RUN | Sat_Jan_29_17:45:23 | – |
| 9150 | WORKER_JOB | d0farm | Worker_D0 | RUN | Sat_Jan_29_18:46:24 | – |
| 9150 | END | d0farm | io_d0sgi | PEND | – | – |
| 9166 | *job* | | | | | |
| 9166 | START_SAM | d0farm | io_d0sgi | RUN | Sun_Jan_30_09:40:51 | – |
| 9166 | WORKER_JOB | d0farm | Worker_D0 | RUN | Sun_Jan_30_09:41:33 | – |
| 9166 | END | d0farm | io_d0sgi | PEND | – | – |
| 9170 | *job* | | | | | |
| 9170 | START_SAM | d0farm | io_d0sgi | RUN | Sun_Jan_30_14:28:38 | – |

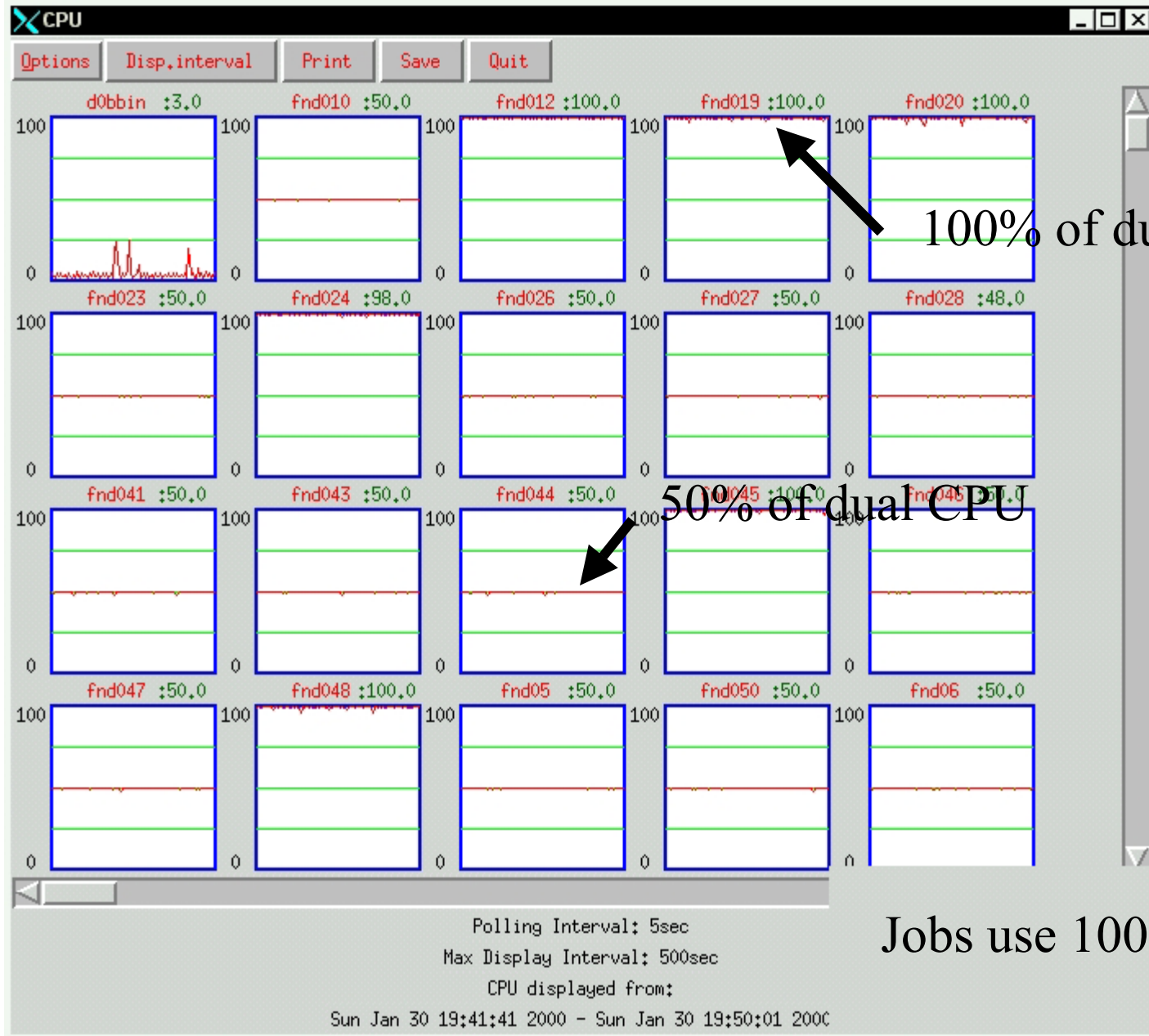## Job Statistics

Job 9150

_____

Step Name: START_SAM
  Host: d0bbin.fnal.gov


  Process Number: 1
  PID    CPU    ACPU   CMD
66443  0      6         /bin/tcsh -f /home/d0farm/aug99/farm_machinery/samtest/start_sam_v5.csh p
66473  6      6           python -u /home/d0farm/aug99/farm_machinery/samtest/start_sam_v5.py pre
162234 0      0             sleep 600
_____

# Farm Batch System Monitor



DO Farms

**CPU** window:
Options | Disp.interval | Print | Save | Quit

| d0bbin :3.0 | fnd010 :50.0 | fnd012 :100.0 | fnd019 :100.0 | fnd020 :100.0 |
| fnd023 :50.0 | fnd024 :98.0 | fnd026 :50.0 | fnd027 :50.0 | fnd028 :48.0 |
| fnd041 :50.0 | fnd043 :50.0 | fnd044 :50.0 | fnd045 :100.0 | fnd046 :100.0 |
| fnd047 :50.0 | fnd048 :100.0 | fnd05 :50.0 | fnd050 :50.0 | fnd06 :50.0 |

100% of dual

50% of dual CPU

Jobs use 100% of CPU

Polling Interval: 5sec
Max Display Interval: 500sec
CPU displayed from:
Sun Jan 30 19:41:41 2000 - Sun Jan 30 19:50:01 2000

File  Edit  View  Go  Communicator  Help

Back | Forward | Reload | Home | Search | Netscape | Print | Security | Shop | Stop

Bookmarks  Location: http://d0ora2.fnal.gov/misweb/cgi/misweb.pl  What's Related

Instant Message | WebMail | Radio | People | Yellow Pages | Download | Calendar | Channels

# SAM Catalog Web Query Interface

## Analyzed Files

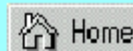| FileName | ConsumerId | Status | ConsumedDate | ProcessId | ProjName | Station | Node |
|---|---|---|---|---|---|---|---|
| sim.pmc02_01.pythia.ztautau_mb1.1av_200evts.276_1151 | 2235 | consumed | 29-jan-00/18:45:04 | 8506 | farmjob.8923 | protofarm | fnd013.fnal, |
| sim.pmc02_01.pythia.ztautau_mb1.1av_200evts.267_1553 | 2235 | consumed | 29-jan-00/18:52:00 | 8507 | farmjob.8923 | protofarm | fnd030.fnal, |
| sim.pmc02_01.pythia.ztautau_mb1.1av_200evts.276_1152 | 2235 | consumed | 29-jan-00/18:53:38 | 8513 | farmjob.8923 | protofarm | fnd031.fnal, |
| sim.pmc02_01.pythia.ztautau_mb1.1av_200evts.267_1552 | 2235 | consumed | 29-jan-00/19:01:19 | 8509 | farmjob.8923 | protofarm | fnd032.fnal, |
| sim.pmc02_01.pythia.ztautau_mb1.1av_200evts.265_1421 | 2235 | consumed | 29-jan-00/19:24:42 | 8508 | farmjob.8923 | protofarm | fnd033.fnal, |

Rows 1 to 5 of the Total 5 found.

Back to: Starting Query Page or [Edit] the SQL query that produced this page.

For help contact sam_support@fnal.gov

**Query to see which input files were processed by a job**

[Home]

**MISWEB Query Interface**

Document: Done

# SAM Catalog Web Query Interface

## Data Files

| FileName | DataTier | CreateDate | RunN |
|---|---|---|---|
| reco.sim.pmc02_01.pythia.ztautau_mb1.1av_200evts.265_1421_8923_4_preco03.05<br>reco.sim.pmc02_01.pythia.ztautau_mb1.1av_200evts.265_1421_8923_4_preco03.05 | reconstructed | 29-JAN-00 | 592 |
| reco.sim.pmc02_01.pythia.ztautau_mb1.1av_200evts.267_1552_8923_5_preco03.05<br>reco.sim.pmc02_01.pythia.ztautau_mb1.1av_200evts.267_1552_8923_5_preco03.05 | reconstructed | 29-JAN-00 | 506 |
| reco.sim.pmc02_01.pythia.ztautau_mb1.1av_200evts.267_1553_8923_3_preco03.05<br>reco.sim.pmc02_01.pythia.ztautau_mb1.1av_200evts.267_1553_8923_3_preco03.05 | reconstructed | 29-JAN-00 | 507 |
| reco.sim.pmc02_01.pythia.ztautau_mb1.1av_200evts.276_1151_8923_1_preco03.05<br>reco.sim.pmc02_01.pythia.ztautau_mb1.1av_200evts.276_1151_8923_1_preco03.05 | reconstructed | 29-JAN-00 | 601 |
| reco.sim.pmc02_01.pythia.ztautau_mb1.1av_200evts.276_1152_8923_2_preco03.05<br>reco.sim.pmc02_01.pythia.ztautau_mb1.1av_200evts.276_1152_8923_2_preco03.05 | reconstructed | 29-JAN-00 | 602 |

Rows 1 to 5 of the Total 5 found.

Back to: Starting Query Page or   Edit   the SQL query that produced this page.

**Check to see if output files were stored properly**

— For help contact sam_support@fnal.gov —   Home
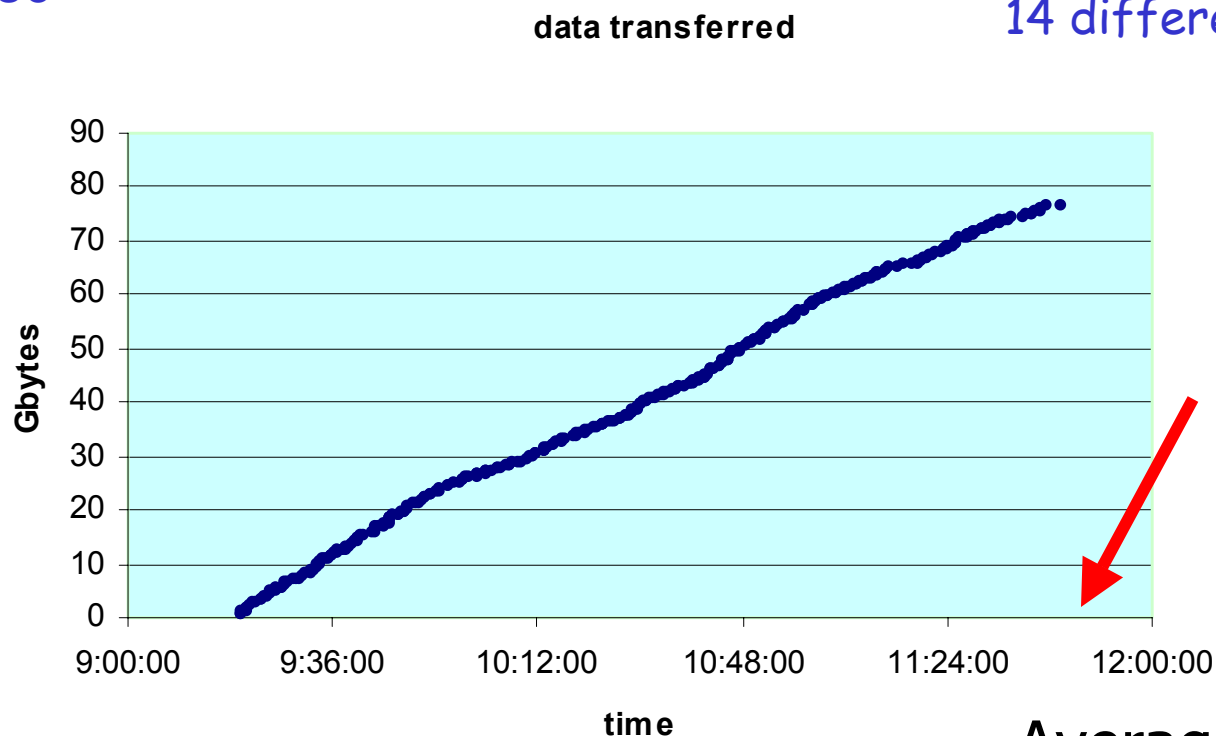
MISWEB Query Interface

Document: Done

# Results of typical farm test

- Create 4 jobs with 25-180 files in each (350 total)
- Submit 4 jobs to the farms using 10-30 workers each (occupy 95/100)
- Process those files through official reconstruction executable
- Files are 200-700 MB Monte Carlo, take 2-10 hours to process.
- 14 tapes read by 5 tape drives (3MB/sec max/drive)
- Output written to I/O node for later dump to tape
- This is almost* equivalent to starting a production 100 processor farm from a cold start.

  *exception is tape drive speed -> 12MB/sec, did not do output to tape

# Data transfer to workers

**Fire up 4 jobs**
**Zee, zmumu, ttbar**
**Qcdpt>80**

322 files
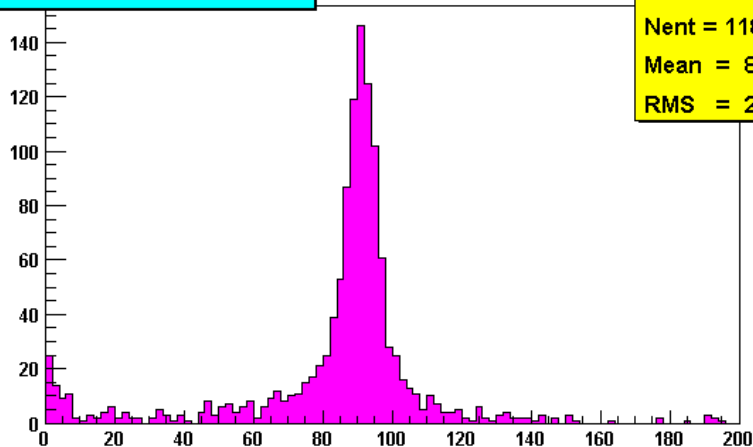95 worker CPU
5 tape drives
14 different tapes

**data transferred**

All files on
workers

Average transfer rate
9.5 MB/sec
Peak ~ 15 MB/sec
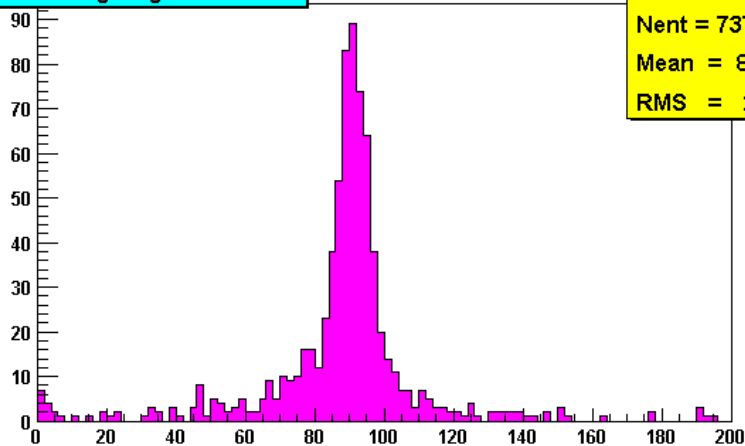
**zmumu_mb1.1av 3800 events**

**DimuonMass**

| hist101 |
|---|
| Nent = 1189 |
| Mean = 84.24 |
| RMS = 27.66 |

**DimuonMass tight-tight**

| hist118 |
|---|
| Nent = 737 |
| Mean = 87.96 |
| RMS = 23.3 |

**Studies done with output of test – muon id validation**

# Results

- Farm system used to debug production executable
  - 8 different causes of crashes found in ~100,000 events
- Improved executable now being used to process 400,000 Monte Carlo events
- FBS/SAM can process data at full rate on 100 processors.
  - System can be scaled by cloning
- Recovery mechanism reasonably robust
- Database can easily track and resubmit failed files.

# Future

- Stop worrying if the farm hardware works, it does
- Use what we learned from first test to optimize submission scripts and file tracking
- Implement file merging on output machine
- Keep testing
- Wait for data

Enstore robot and PC mover nodes

SUN 4500

SAM DB and PC servers

Cisko 6509 Switch

SGI 02000

50 Dual PC worker nodes

DO Farms